

Open Research Online

The Open University's repository of research publications and other research outputs

Factors in human recognition of timbre lexicons generated by data clustering

Conference or Workshop Item

How to cite:

Roma, Gerard; Xambó, Anna; Herrera, Perfecto and Laney, Robin (2012). Factors in human recognition of timbre lexicons generated by data clustering. In: 9th Sound and Music Computing Conference (SMC 2012), 11-14 Jul 2012, Copenhagen, Denmark.

For guidance on citations see [FAQs](#).

© 2012 Gerard Roma et al.

Version: Not Set

Link(s) to article on publisher's website:
<http://smc2012.smcnetwork.org/>

Copyright and Moral Rights for the articles on this site are retained by the individual authors and/or other copyright owners. For more information on Open Research Online's data [policy](#) on reuse of materials please consult the policies page.

oro.open.ac.uk

Factors in human recognition of timbre lexicons generated by data clustering

Gerard Roma

Music Technology Group
Universitat Pompeu Fabra
gerard.roma@upf.edu

Anna Xambó

Music Computing Lab
Open University
a.xambo@open.ac.uk

Perfecto Herrera

Music Technology Group
Universitat Pompeu Fabra
perfecto.herrera@upf.edu

Robin Laney

Music Computing Lab
Open University
r.c.laney@open.ac.uk

ABSTRACT

Since the development of sound recording technologies, the palette of sound timbres available for music creation was extended way beyond traditional musical instruments. The organization and categorization of timbre has been a common endeavor. The availability of large databases of sound clips provides an opportunity for obtaining data-driven timbre categorizations via content-based clustering. In this article we describe an experiment aimed at understanding what factors influence the process of learning a given clustering of sound samples. We clustered a large database of short sound clips, and analyzed the success of participants in assigning sounds to the “correct” clusters after listening to a few examples of each. The results of the experiment suggest a number of relevant factors related both to the strategies followed by users and to the quality measures of the clustering solution, which can guide the design of creative applications based on audio clip clustering.

1. INTRODUCTION

Web-based sound databases provide an interesting resource for facilitating music and multimedia creation. The use of sound samples makes it easier for people without specific musical or technical training to create audio content. Web-based content management applications are allowing internet users to contribute samples as basic building blocks. Some examples include Freesound.org, which has now around 150k sounds released under Creative Commons licenses, or Looperman.com, where one can download around 32k royalty-free loops and samples. However, the lack of a traditional editorial process creates new challenges for the organization and access to these samples. Automatic data clustering has been used in many domains to provide an intuitive interface to large data collections. Clustering techniques can find existing partitions in data

that is not labeled or previously classified. Hence, they are a common choice as a means for exploration of multimedia databases.

Yet, the question of how meaningful for users are the clusters found in data is often overlooked. Clustering algorithms are usually evaluated using either external or internal criteria [1]. External criteria compare partitions obtained in the clustering process to previously defined categories. These criteria allow to evaluate algorithms for known partitions, which give a hint about how well they will perform with unlabeled data. Internal criteria measure the quality of the solution in terms of the distance metrics defined for the data. In many cases, labels are not available for the data (which is the reason for using a clustering algorithm). Internal clustering quality measures do not require labels, but they do not provide the whole picture of how useful the solutions will be for humans. Other factors, such as users backgrounds and their understanding of the data as presented through the interface, influence the usefulness of applications based on data clustering.

In this article we analyze which factors influence the success of users in recognizing the groupings produced by an automatic sound sample clustering algorithm. We analyze recognition rate and how it is related both to internal quality measures of the clustering solution and to other factors related with each individual user, such as the strategy followed for learning the clusters. We investigate human factors through video analysis of an interactive exercise consisting on listening to examples of each cluster and assigning new sounds to one of the partitions. We conclude giving some design considerations derived from our analysis.

2. TIMBRE LEXICONS FOR MUSIC CREATION

The exploration of the timbre space opened by electronic recordings started in the field of contemporary music with the analysis of Pierre Schaeffer [2], who proposed a foundation for a new musical theory based on acoustic criteria. Several decades earlier, Luigi Russolo had described a classification of noises as the basis of a futurist music [3]. With the evolution of sound media technologies, composers have continued experimenting with their own sound alphabets. For example Roads [4] and Lehdal [5]

have studied the use of timbre alphabets in the context of grammars. Smalley [6] proposed spectro-morphology as the analysis of sound categorization.

With the popularization of digital electronic instruments, such as wavetable synthesizers or samplers, the concept of a palette of pre-set waveforms became possibly the most common way to deal with timbre in electronic music creation. This way of working is still widely popular in computer-based music making, and for years specialized companies have sold sample CDs for use in hardware/software instruments. The distribution of waveforms in these instruments has traditionally been tied to technical issues, and often to the cultural impact of some instruments, such as e.g. the Roland TR-909 synthetic drum kit.

Given the amounts of data available in the internet era, sound categorizations can be obtained from large databases through computational means. This approach makes it possible to obtain lexicons based on psychoacoustic features (and not only on technical aspects), that can be used by the general public, while lexicons designed by composers are very often particular to their work. Casey [7] proposed an approach based on Hidden Markov Models for indexing internet audio using spectral prototypes. However, no evaluation was done with respect to whether such prototypes would be meaningful for users. Since clustering algorithms can always provide partitions of data, it is important to understand what factors determine their usability in creative applications.

3. CONTENT-BASED CLUSTERING OF AUDIO CLIPS

One simple way of obtaining categorizations of sounds is through data clustering. Clustering is an unsupervised process that finds existing divisions in data. Given a sufficiently large database, it should be possible to obtain a rich palette of different types of sounds. While obtaining an optimal method for clustering audio clips is not the focus of this paper, we briefly describe the main elements of the used method.

3.1 Related work

Clustering algorithms have been mainly used for visual exploration of audio databases [8]. Several works have described the use of Self Organizing Maps (SOM) which provide a graphical representation, including applications to drum sample collections [9], music files [10] and sound effects [11]. Most of these works focus on visual exploration, and do not evaluate whether the divisions presented by the SOM are understood by users. Our focus is different in that we analyze the grouping produced by the algorithm. Our goal is to understand whether this approach could provide support for interacting with audio databases beyond the initial exploration step, by allowing the user to become familiar with a lexicon of sound categories.

3.2 Audio features

The process of clustering is based on a measure of similarity between the audio clips. The most common ap-

proach is to compute this similarity from some distance metric between feature vectors extracted from the audio signal. In the SOM-based clustering literature there are several works that use the amount of energy of the signal in different frequency bands of the bark scale [12]. The result is a quantized spectrogram that approximates human perception. We use a similar approach but based on the gammatone filter bank [13]. Our implementation is based in the frequency domain approximation proposed by Ellis [14], which allows quickly analyzing massive collections of sound clips.

We compute the gammatone features for successive windows of the audio signal, using a window size of 23ms and a hop size of 11ms. In order to represent the whole audio clip, we then compute the mean and variance of each band as well as its first and second derivatives. In order to obtain more temporal information, we extract the modulation spectrum, as suggested in [15], by computing the magnitude spectrum of the temporal evolution of each band, and adding the energies of that spectrum to three bands (1–2Hz, 3–15Hz, 20–43Hz, we ignore the DC component at 0Hz).

This results in vectors of large dimensionality (162). In order to work with such high-dimensional data, we use cosine distance to determine the similarity between audio clips, and an algorithm based on k-nearest neighbor graph as described in the next section.

3.3 Clustering algorithm

In order to obtain the clusters, we use sounds from Freesound.org, a collaborative database of sound clips released under Creative Commons licenses. This allows us to work with a large collection of samples, but also bears problems of noise and uneven feature densities that are common when dealing with data from the web. Many algorithms have been developed, beyond traditional approaches, to deal with this kind of data. We use a partial implementation of the Chameleon algorithm available in the CLUTO [16] package. This algorithm starts by computing a k-nearest neighbor graph from the data points, which is then partitioned using a min-cut algorithm [17]. Using the k-nearest neighbor graph allows us to find clusters of different densities. For example, one user may upload many sounds from the same source recorded in the same conditions, which will produce a very tight cluster. On the other hand, sounds of some particular instrument may be more difficult to find and may be recorded in different conditions. Distances between them are generally greater, forming much sparser clusters.

Internal quality measures are often used to evaluate the obtained clustering solution. One common approach is to analyze the similarities within each cluster to understand how compact it is, and the similarities between points of each cluster and all points in the others to measure how well separated are the different clusters in the solution. As mentioned, we use cosine similarity, i.e.:

$$S_{ij} = \frac{d_i d_j}{|d_i| |d_j|} \quad (1)$$

where d_i, d_j are the feature vectors extracted for sounds i

and j . Then given the cluster C_n (the n th set of feature vectors in the clustering solutions, where $n \in (1..N)$ if N clusters were found), we can define:

$$C_{sim}(C_n) = \frac{\sum_{(d_i, d_j) \in C_n} S_{ij}}{|C_n|} \quad (2)$$

$$C_{imax}(C_n) = \frac{\max_{(d_i, d_j) \in C_n} S_{ij}}{|C_n|} \quad (3)$$

$$C_{imin}(C_n) = \frac{\min_{(d_i, d_j) \in C_n} S_{ij}}{|C_n|} \quad (4)$$

respectively the mean, minimum and maximum similarity between points within cluster C_n . The mean similarity gives an indication of how compact is the cluster overall. The minimum similarity corresponds to the maximum distance between two points, which can indicate the presence of outliers. The maximum similarity in the cluster can be used as a hint of the maximum density inside the cluster, i.e. if the maximum similarity is small then the cluster may be sparse, whereas a dense cluster with some outliers may have a smaller average similarity but still a high maximum.

Analogously, $C_{esim}(C_n)$, $C_{emin}(C_n)$ and $C_{emax}(C_n)$ are based on external similarities, i.e. similarities between points in C_n and points in all other clusters.

4. USER EXPERIMENT

4.1 Motivation

As described in the previous sections, the motivation of this study is to understand what factors influence the usability of a given sound clip clustering approach. Our main assumption is that this usability is determined by how well the user is able to learn the groupings created by the clustering algorithm. We tested this assumption by analyzing the users in the tasks of learning the clusters found by the algorithm and assigning unlabeled sounds to their cluster. Findings of the experiment should be valuable for choosing the appropriate clustering approach as much as in the design of interfaces for interacting with large sound clip databases, mainly for creative applications.

4.2 Database sample

Our aim is analyzing the viability of a clustering scheme for interacting with large databases. In most applications we expect that a fast learning process is necessary for the user to retain the interest (we focus on applications where the user does not need a specific background or training). Thus, in order to test a realistic use case, we sampled the underlying clustered dataset for the experiment.

In order to obtain the test dataset, we first clustered a large database of 10k sounds obtained from Freesound.org, all shorter than one second in order to avoid sounds with many acoustic events and timbre variation. Determining an optimal number of clusters is a non-trivial issue that was not the focus in this work. We chose a number of clusters that in a subjective evaluation gave consistent clusters while allowing a manageable size for the lexicon (in the order of

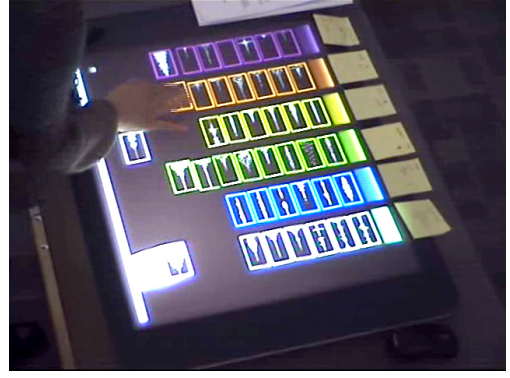


Figure 1. Picture of the prototype.

e.g. the size of the latin alphabet or the number of keys in a keyboard). In the process it became clear that a larger number would give smaller and more consistent clusters. Yet, in real world applications we can not expect the user to learn hundreds of sound categories. We ran our algorithm to produce 40 clusters. Of these, we discarded clusters with less than 50 instances and chose a random sample of 6 clusters for each user. This number seems well aligned with acceptable cognitive load in short term memory [18]. Of these clusters, we randomly chose 6 example sounds for each and again chose 20 random sounds from the pool of chosen clusters as test for that user.

4.3 Prototype

For the experiment, we implemented a simple prototype on a Microsoft Surface multi-touch table, using Adobe Air. The interface showed 6 colored rectangles representing the different clusters, each of which could be unfolded to visualize and play the sounds. Hence, the only potential visual cues with respect to the clusters were the sound waveforms of the examples. These images are the same that are used for sound lists in Freesound.org (except that we removed the color). Test examples were showed as a pile of black and white waveform objects resembling a deck of cards. Dragging each card to the vertical area below the color shape corresponding to each cluster colored the object to the cluster color, which signaled that the sound was assigned to that cluster.

4.4 Experiment protocol

The experiment was divided in two main tasks. In task 1, participants were asked to listen to the sound examples of each cluster. In addition, they were asked to annotate any words or tags needed to identify and remember each cluster in sticky paper notes that were attached to the table above each cluster area. This allowed us to analyze how the users understood the clusters. In task 2, participants were asked to classify the stack of test sounds into the clusters by dragging them to the appropriate area.

The use of a large multi-touch table provided a more embodied interaction that allowed us to observe and analyze participants movements and strategies. Video was recorded with two cameras positioned non-intrusively:

general view and close-up view. Video output from the device was also captured for complementing the analysis.

Finally, a questionnaire was filled in with basic demographic information, as well as some questions about their own confidence in the performed task and the criteria they followed for the classification.

4.5 Participants

The study took part in a computing department and most of the participants were familiar with computers, although not necessarily with music or audio-related topics. In total there were 14 participants (9 males, 5 females) with ages from 21 to 50 and a diversity of nationalities and cultural backgrounds. Most had some kind of musical training: 4 reported no training at all, 4 some degree of music training and 6 of them more than five years. With respect to familiarity with electronic music creation, 8 of them reported no previous experience, 5 of them had some familiarity, and one of them was a regular user of electronic music creation tools.

4.6 Data analysis

After the experiment, we had several sources of data for analysis. The prototype logged the classification choices performed by users, and kept the information about the clusters that were presented. On the other hand, the footage from the cameras was used for video analysis, which is a common tool in human-computer interaction studies [19]. This technique allows for a great deal of detail if compared to more traditional HCI methods, but requires time and a clear focus with respect to the problem at hand. We followed observations made during the experiment and an initial overview of the video to define our target variables (described in the next section). A coding scheme for annotating the video was then defined by two of the authors. The *Elan*¹ software was used for the video analysis. The main advantage of this program is that it allows hierarchical specification of the code for annotating the video. We used the annotations, along with the data from the questionnaire, for qualitative and quantitative analysis of user-related factors. Finally, we used the data from the clustering and the logged results for quantitative analysis of factors related with the clusters.

5. RESULTS AND DISCUSSION

We analyzed both qualitative and quantitative aspects of the experiment in order to understand which factors could determine the correct assignment of sounds to their own cluster as computed by the clustering algorithm. This assignment was encoded as a binary variable for each of the tested sounds. Each participant was given 20 sounds, so in total there were 280 assignments. We aggregated the results in order to analyze the data from the point of view of the user and from the point of view of the clustering algorithm. In the first case, the target variable became the fraction of correctly assigned sounds by each user, and in the

second, the fraction of correctly assigned sounds for each cluster. Table 1 shows a summary of all variables resulting from the experiment.

5.1 User level

To understand human factors related to the proposed task, we did a qualitative analysis of the responses to the questionnaire, as well as an analysis of the video footage. The fraction of successfully assigned sounds of each user oscillated around 40% ($mean = 0.44$, $sd = 0.5$).

We found several recurrent themes in the analysis of the questionnaire responses referring to the criteria used to classify the sounds. They are summarized in Table 2. Most popular criteria could be classified as “Sound sources” and “Sound properties”.

In the video analysis we observed some relevant aspects related with the behavior of participants. Our main observation was that participants followed different strategies in both tasks. In task 1, there were 3 participants who started by pre-listening to all or most of the sounds on the table before starting. This group devoted more time in the two tasks, and one of them scored the best result. Another difference was found between those participants who explored neighbor clusters before labeling a given one, and those who did not. In task 2, we observed that some participants tended to rely on their written notes, while others went back to listening to the examples when deciding to which cluster they would assign the sound. All but one of the participants who correctly assigned more than 50% of the sounds followed this strategy of constantly comparing the test sound with examples.

We confirmed the significance of these differences through two-sample t-tests over the three binary variables: users who followed strategy S_{t1s1} (pre-listening to the whole table) took longer ($p < 1e^{-15}$) and did better ($p < 1e^{-04}$), as well as users who followed strategy S_{t1s2} (looking at contiguous clusters) ($p < 1e^{-07}$). The result for S_{t2s2} showed that the group that kept listening to sounds performed better ($p < 0.02$). We further counted the number of times the example sounds and test sounds were played for the first 10 test sounds (after that, users tended to classify without playing the examples again). The overall count of plays of the examples in task 2 for each user correlates with the recognition rate ($r = 0.64$), while the number of plays of the test sound appeared to be barely correlated ($r = 0.122$). In all, we were able to extract more significant variables from the learning phase (task 1), and the relevant outcomes in the classification phase still seemed to refer to learning strategies, which reflects the importance of this step.

We focused with some more detail on the three participants that scored best. All of them referenced concepts related to sound sources as their criteria in the questionnaire. Their notes taken during task 1 tended to be more consistent and easier to compare. During each classification task, they tended to maintain attention until they located the target cluster. One technique that was particular to these users was fast pre-listening of several example sounds in a cluster, which produced a quick audio summary of that cluster.

¹ <http://www.lat-mpi.eu/tools/elan>, developed at the Max Planck Institute for Psycholinguistics, Nijmegen, The Netherlands [20].

Cluster level	
C_{size}	Cluster size ($ C_n $)
C_{isim}	Average internal similarity
C_{imax}	Maximum internal similarity
C_{imin}	Minimum internal similarity
C_{esim}	Average external similarity
C_{emax}	Maximum external similarity
C_{emin}	Minimum external similarity
User level	
S_{t1s1}	pre-listening of all table
S_{t1s2}	listening to neighbor clusters
S_{t2s1}	listening to examples again

Table 1. Variables obtained from the experiment.

Themes	Classification criteria
Source of the sound (16)	<ul style="list-style-type: none"> - instruments (5): e.g., <i>bass</i> (2), <i>drums</i> (3) - onomatopoeias (1): e.g., <i>click click</i> (1) - <i>speech/vocal</i> (2) - <i>non-speech</i> (1) - physical source (2): <i>physical phenomena</i> (1) / <i>physical objects</i> (1) - electronic sounds (4): e.g., <i>synthesizer banks</i> (1), <i>base/moog</i> (1), <i>sound effect</i> (1) / <i>synthetic or futuristic sounds</i> (1) - <i>everyday sounds</i> (1)
Sound properties (10)	<ul style="list-style-type: none"> - <i>pitch/brightness</i> (2) / <i>tone</i> (2) - <i>length</i> (3) - dichotomies (2): e.g., <i>hard/soft</i> (1), <i>loud/low</i> (1) - <i>sound envelope</i> (1)
Experiential (5)	<ul style="list-style-type: none"> - <i>general feel</i> (1) - <i>instinct/intuition</i> (2) - <i>stories</i> (1) - <i>music experience/knowledge</i> (1)
Sound and/or visual similarity (5)	<ul style="list-style-type: none"> - sound similarity (2): e.g., <i>similar sound features</i> (2) - <i>similarity</i> (1) - visual similarity (2): e.g., <i>waveform</i> (1) / <i>waveform icons</i> (1)
Sound description (3)	- categories/tags (3): e.g., <i>my own postits</i> (1), <i>categories</i> (2)
Overall sound (1)	- <i>mainly the sound</i> (1)

Table 2. Themes and classification criteria followed by users extracted from the questionnaire.

ter. After some iterations, the different clusters had been learned and pre-listening was no longer necessary.

5.2 Cluster level

When looking at the results aggregated from the point of view of clusters, recognition rate was similar but with less variation ($mean = 0.48$, $sd = 0.2$). In order to understand the importance of different measures of cluster quality (outlined in section 3.3), we built a multiple linear regression model using these measures as independent variables and the recognition rate as dependent variable. One common problem with linear regression is multicollinearity due to correlation of the independent variables, which can give misleading results. We checked the Variable Inflation Factor (VIF), for controlling multicollinearity problems, and ensured that it was below 10, which is the usually recommended threshold [21]. This forced us to re-

move the C_{imin} and C_{emin} variables, which represent the maximum internal/external distance (minimum similarity) and thus are highly related to the corresponding mean variables. We also removed C_{emax} , as it didn't make any significant contribution to the model. Table 3 shows the coefficients of the model. Perhaps surprisingly, C_{imax} , the maximum similarity within the cluster, has the highest significant impact over the recognition rate, much higher than the average similarity. This means that clusters with high maximum similarity were easy to learn while clusters with low maximum similarity were difficult. In relation with the lower weight of the mean similarity, this suggests that clusters containing high-density areas allow an easier association of acoustic features with a single label, while sparse clusters, where the closest two points are not particularly close (even if the average similarity stays high) should be avoided. The rest of coefficients are more or less

predictable. The size of the cluster has a small but significant negative effect, which suggests that smaller clusters are to be preferred.

Variable	Estimate	Std. Error	t value	Pr(> t)
C_{imean}	1.1229	0.6682	1.68	0.1048
C_{imax}	22.5082	9.6593	2.33	0.0278
C_{size}	-0.0009	0.0004	-2.19	0.0378
C_{emean}	-2.1158	0.8856	-2.39	0.0244
Adjusted R^2	0.3666			

Table 3. Regression analysis for cluster-level variables.

6. CONCLUSIONS

Websites that focus on user-contributed sound clips are creating an interesting resource for creative applications related to audio. Perhaps the main challenge for interacting with these collections is the noise (or even worse, silence) in the description and organization of the sounds. Providing an appropriate textual description is a tedious, and often overlooked, part of the process. In this context, content-based data clustering has the potential of automatically creating meaningful partitions. Our experiment provided valuable insights in this respect. In particular, it highlighted the differences between the proposed use case of learning a lexicon of sound object categories and that of visually exploring a database, which is a common application of data clustering.

From the point of view of users, the result of our experiment stressed the importance of the learning phase, where most significant differences between users were observed. This suggests that interfaces for applications based on clustering could make use of a specialized interface for learning the clusters. Also, interfaces should make it easy to compare the sounds of different clusters. Finally, it seems that references to sound sources as well as key acoustic properties in each cluster are common labels that users associate with the partitions.

From the point of view of the clustering algorithm, the experiment suggests that algorithms should concentrate on (potentially small) areas of high density, perhaps discarding a certain number of outliers if the database is large enough. Still, this is not a trivial issue with heterogeneous data such as the sounds in Freesound.org. We are currently working on a suitable approach based on the insights gained in the present work.

Acknowledgments

The authors would like to thank Yvonne Rogers for providing us access to the multi-touch table used in the test, Chris Dobbyn for his insightful comments, and all the participants for their time and feedback.

7. REFERENCES

- [1] A. K. Jain and R. C. Dubes, *Algorithms for clustering data*. Upper Saddle River, NJ, USA: Prentice-Hall, Inc., 1988.
- [2] P. Schaeffer, *Traité des objets Musicaux*. Paris: Seuil, 1966.
- [3] L. Russolo, *The Art of Noises*. New York: Pendragon Press, 1986.
- [4] C. Roads, “Composing grammars,” in *Proc. Int. Computer Music Conf. (ICMC)*, 1978.
- [5] F. Lerdahl, “Timbral hierarchies,” *Contemporary Music Review*, vol. 2, no. 1, pp. 135–160, 1987.
- [6] D. Smalley, “Spectromorphology: explaining sound-shapes,” *Organised Sound*, vol. 2, no. 2, pp. 107–126, July 1997.
- [7] M. Casey, “Acoustic lexemes for organizing internet audio,” *Contemporary Music Review*, vol. 24, no. 6, pp. 489–508, December 2005.
- [8] M. Cooper, J. Foote, E. Pampalk, and G. Tzanetakis, “Visualization in audio-based music information retrieval,” *Computer Music J.*, vol. 30, pp. 42–62, June 2006.
- [9] E. Pampalk, P. Hlavac, and P. Herrera, “Hierarchical organization and visualization of drum sample libraries,” in *Proc. Int. Conf. Digital Audio Effects (DAFx)*, 2004, pp. 378–383.
- [10] E. Pampalk, A. Rauber, and D. Merkl, “Content-based organization and visualization of music archives,” in *Proc. ACM Int. Conf. Multimedia (ACM-MM)*, 2002, pp. 570–579.
- [11] E. Brazil, M. Fernstroem, G. Tzanetakis, and P. Cook, “Enhancing sonic browsing using audio information retrieval,” in *Proc. Int. Conf. Auditory Display (ICAD)*, 2002, pp. 132–135.
- [12] E. Zwicker and H. Fastl, *Psychoacoustics, Facts and Models*. Berlin: Springer-Verlag, 1990.
- [13] R. F. Lyon, A. G. Katsiamis, and E. M. Drakakis, “History and future of auditory filter models,” in *Proc. IEEE Int. Symp. Circuits and Systems (ISCAS)*, 2010, pp. 3809–3812.
- [14] D. P. W. Ellis, “Gammatone-like spectrograms,” 2009. [Online]. Available: <http://www.ee.columbia.edu/~dpwe/resources/matlab/gammatonegram/>
- [15] M. F. McKinney and J. Breebaart, “Features for audio and music classification,” in *Proc. Int. Symp. Music Information Retrieval (ISMIR)*, 2003, pp. 151–158.
- [16] G. Karypis, “CLUTO - a clustering toolkit,” University of Minnesota, Department of Computer Science, Tech. Rep. #02-017, Apr. 2002.
- [17] G. Karypis, E.-H. Han, and V. Kumar, “Chameleon: hierarchical clustering using dynamic modeling,” *Computer*, vol. 32, no. 8, pp. 68–75, Aug. 1999.

- [18] G. Miller, "The magical number seven, plus or minus two: Some limits on our capacity for processing information," *The Psychological Review*, vol. 63, pp. 81–97, 1956.
- [19] B. Jordan and A. Henderson, "Interaction analysis: Foundations and practice," *The J. of the Learning Sciences*, vol. 4, no. 1, pp. 39–103, 1995.
- [20] H. Sloetjes and P. Wittenburg, "Annotation by category: ELAN and ISO DCR," in *Proc. Int. Conf. Language Resources and Evaluation (LREC)*, 2008, pp. 816–820.
- [21] J. Hair, W. Black, B. Babin, and R. Anderson, *Multivariate data analysis: a global perspective*, 7th ed. London: Pearson Education, 2010.